

A Generalized Representer Theorem

Bernhard Schölkopf^{1,2,3}, Ralf Herbrich^{1,2}, and Alex J. Smola¹

¹ Australian National University, Department of Engineering, Canberra ACT 0200, Australia

`bs@conclu.de`, `rherb@microsoft.com`, `alex.smola@anu.edu.au`

² Microsoft Research Ltd., 1 Guildhall Street, Cambridge, UK

³ New address: Biowulf Technologies, Floor 9, 305 Broadway, New York, NY 10007, USA

Abstract. Wahba's classical representer theorem states that the solutions of certain risk minimization problems involving an empirical risk term and a quadratic regularizer can be written as expansions in terms of the training examples. We generalize the theorem to a larger class of regularizers and empirical risk terms, and give a self-contained proof utilizing the feature space associated with a kernel. The result shows that a wide range of problems have optimal solutions that live in the finite dimensional span of the training examples mapped into feature space, thus enabling us to carry out kernel algorithms independent of the (potentially infinite) dimensionality of the feature space.

1 Introduction

Following the development of support vector (SV) machines [23], positive definite kernels have recently attracted considerable attention in the machine learning community. It turns out that a number of results that have now become popular were already known in the approximation theory community, as witnessed by the work of Wahba [24]. The present work brings together tools from both areas. This allows us to formulate a generalized version of a classical theorem from the latter field, and to give a new and simplified proof for it, using the geometrical view of kernel function classes as corresponding to vectors in linear feature spaces.

The paper is organized as follows. In the present first section, we review some basic concepts. The two subsequent sections contain our main result, some examples and a short discussion.

1.1 Positive Definite Kernels

The question under which conditions kernels correspond to dot products in linear spaces has been brought to the attention of the machine learning community by Vapnik and coworkers [1,5,23]. In functional analysis, the same problem has been studied under the heading of Hilbert space representations of kernels. A good monograph on the functional analytic theory of kernels is [4]. Most of the material in the present introductory section is taken from that work. Readers familiar with the basics of kernels can skip over the remainder of it.

Suppose we are given empirical data

$$(x_1, y_1), \dots, (x_m, y_m) \in \mathcal{X} \times \mathbf{R}. \quad (1)$$

Here, the target values y_i live in \mathbf{R} , and the patterns x_i are taken from a domain \mathcal{X} . The only thing we assume about \mathcal{X} is that is a nonempty set. In order to study the problem of learning, we need additional structure. In kernel methods, this is provided by a similarity measure

$$k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbf{R}, \quad (x, x') \mapsto k(x, x'). \quad (2)$$

The function k is called a kernel [20]. The term stems from the first use of this type of function in the study of integral operators, where a function k giving rise to an operator T_k via

$$(T_k f)(x) = \int_{\mathcal{X}} k(x, x') f(x') \, dx' \quad (3)$$

is called the kernel of T_k . Note that we will state most results for the more general case of complex-valued kernels;¹ they specialize to the real-valued case in a straightforward manner. Below, unless stated otherwise, indices i and j will be understood to run over the training set, i.e. $i, j = 1, \dots, m$.

Definition 1 (Gram matrix). *Given a kernel k and patterns $x_1, \dots, x_m \in \mathcal{X}$, the $m \times m$ matrix*

$$K := (k(x_i, x_j))_{ij} \quad (4)$$

is called the Gram matrix of k with respect to x_1, \dots, x_m .

Definition 2 (Positive definite matrix). *An $m \times m$ matrix K over the complex numbers \mathbf{C} satisfying*

$$\sum_{i=1}^m \sum_{j=1}^m c_i \bar{c}_j K_{ij} \geq 0 \quad (5)$$

for all $c_1, \dots, c_m \in \mathbf{C}$ is called positive definite.

Definition 3 (Positive definite kernel). *Let \mathcal{X} be a nonempty set. A function $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbf{C}$ which for all $m \in \mathbf{N}$, $x_1, \dots, x_m \in \mathcal{X}$ gives rise to a positive definite Gram matrix is called a positive definite (pd) kernel.²*

Real-valued kernels are contained in the above definition as a special case. However, it is not sufficient to require that (5) hold for real coefficients c_i . If we want to get away with real coefficients only, we additionally have to require that the

¹ We use the notation \bar{c} to denote the complex conjugate of c .

² One might argue that the term *positive definite kernel* is slightly misleading. In matrix theory, the term *definite* is sometimes used to denote the case where equality in (5) only occurs if $c_1 = \dots = c_m = 0$. Simply using the term *positive kernel*, on the other hand, could be confused with a kernel whose *values* are positive.

kernel be symmetric. The complex case is slightly more elegant; in that case, (5) can be shown to imply symmetry, i.e. $k(x_i, x_j) = \overline{k(x_j, x_i)}$.

Positive definite kernels can be regarded as generalized dot products. Indeed, any dot product is a pd kernel; however, linearity does not carry over from dot products to general pd kernels. Another property of dot products, the Cauchy-Schwarz inequality, does have a natural generalization: if k is a positive definite kernel, and $x_1, x_2 \in \mathcal{X}$, then

$$|k(x_1, x_2)|^2 \leq k(x_1, x_1) \cdot k(x_2, x_2). \quad (6)$$

1.2 ... and Associated Feature Spaces

We define a map from \mathcal{X} into the space of functions mapping \mathcal{X} into \mathbf{C} , denoted as $\mathbf{C}^{\mathcal{X}}$, via [4]

$$\phi : \mathcal{X} \rightarrow \mathbf{C}^{\mathcal{X}}, \quad x \mapsto k(\cdot, x). \quad (7)$$

Here, $\phi(x) = k(\cdot, x)$ denotes the function that assigns the value $k(x', x)$ to $x' \in \mathcal{X}$. Applying ϕ to x amounts to representing it by its similarity to *all* other points in the input domain \mathcal{X} . This seems a very rich representation, but it turns out that the kernel allows the computation of a dot product in that representation.

We shall now construct a dot product space containing the images of the input patterns under ϕ . To this end, we first need to endow it with the linear structure of a vector space. This is done by forming linear combinations of the form

$$f(\cdot) = \sum_{i=1}^m \alpha_i k(\cdot, x_i), \quad g(\cdot) = \sum_{j=1}^{m'} \beta_j k(\cdot, x'_j). \quad (8)$$

Here, $m, m' \in \mathbf{N}$, $\alpha_1, \dots, \alpha_m, \beta_1, \dots, \beta_{m'} \in \mathbf{C}$ and $x_1, \dots, x_m, x'_1, \dots, x'_{m'} \in \mathcal{X}$ are arbitrary. A dot product between f and g can be constructed as

$$\langle f, g \rangle := \sum_{i=1}^m \sum_{j=1}^{m'} \bar{\alpha}_i \beta_j k(x_i, x'_j). \quad (9)$$

To see that this is well-defined, although it explicitly contains the expansion coefficients (which need not be unique), note that $\langle f, g \rangle = \sum_{j=1}^{m'} \beta_j \overline{f(x'_j)}$, using $k(x'_j, x_i) = \overline{k(x_i, x'_j)}$. The latter, however, does not depend on the particular expansion of f . Similarly, for g , note that $\langle f, g \rangle = \sum_i \bar{\alpha}_i g(x_i)$. This also shows that $\langle \cdot, \cdot \rangle$ is anti-linear in the first argument and linear in the second one. It is symmetric, since $\langle f, g \rangle = \overline{\langle g, f \rangle}$. Moreover, given functions f_1, \dots, f_n , and coefficients $\gamma_1, \dots, \gamma_n \in \mathbf{C}$, we have

$$\sum_{i=1}^n \sum_{j=1}^n \bar{\gamma}_i \gamma_j \langle f_i, f_j \rangle = \left\langle \sum_{i=1}^n \gamma_i f_i, \sum_{j=1}^n \gamma_j f_j \right\rangle \geq 0, \quad (10)$$

hence $\langle \cdot, \cdot \rangle$ is actually a pd kernel on our function space.

For the last step in proving that it even is a dot product, one uses the following interesting property of ϕ , which follows directly from the definition: for all functions (8),

$$\langle k(\cdot, x), f \rangle = f(x), \quad (11)$$

i.e., k is the *representer of evaluation*. In particular, $\langle k(\cdot, x), k(\cdot, x') \rangle = k(x, x')$, the *reproducing kernel property* [2,4,24], hence (cf. (7)) we indeed have $k(x, x') = \langle \phi(x), \phi(x') \rangle$. Moreover, by (11) and (6) we have

$$|f(x)|^2 = |\langle k(\cdot, x), f \rangle|^2 \leq k(x, x) \cdot \langle f, f \rangle. \quad (12)$$

Therefore, $\langle f, f \rangle = 0$ implies $f = 0$, which is the last property that was left to prove in order to establish that $\langle \cdot, \cdot \rangle$ is a dot product.

One can complete the space of functions (8) in the norm corresponding to the dot product, i.e., add the limit points of sequences that are convergent in that norm, and thus gets a Hilbert space H_k , usually called a *reproducing kernel Hilbert space (RKHS)*. The case of real-valued kernels is included in the above; in that case, H_k can be chosen as a real Hilbert space.

2 The Representer Theorem

As a consequence of the last section, one of the crucial properties of kernels is that even if the input domain \mathcal{X} is only a set, we can nevertheless think of the pair (\mathcal{X}, k) as a (subset of a) Hilbert space. From a mathematical point of view, this is attractive, since we can thus study various data structures (e.g., strings over discrete alphabets [26,13,18]) in Hilbert spaces, whose theory is very well developed. From a practical point of view, however, we now face the problem that for many popular kernels, the Hilbert space is known to be infinite-dimensional [24]. When training a learning machine, however, we do not normally want to solve an optimization problem in an infinite-dimensional space.

This is where the main result of this paper will be useful, showing that a large class of optimization problems with RKHS regularizers have solutions that can be expressed as kernel expansions in terms of the training data. These optimization problems are of great interest for learning theory, both since they comprise a number of useful algorithms as special cases and since their statistical performance can be analyzed with tools of learning theory (see [23,3], and, more specifically dealing with regularized risk functionals, [6]).

Theorem 1 (Nonparametric Representer Theorem). *Suppose we are given a nonempty set \mathcal{X} , a positive definite real-valued kernel k on $\mathcal{X} \times \mathcal{X}$, a training sample $(x_1, y_1), \dots, (x_m, y_m) \in \mathcal{X} \times \mathbf{R}$, a strictly monotonically increasing real-valued function g on $[0, \infty[$, an arbitrary cost function $c : (\mathcal{X} \times \mathbf{R}^2)^m \rightarrow \mathbf{R} \cup \{\infty\}$, and a class of functions*

$$\mathcal{F} = \left\{ f \in \mathbf{R}^{\mathcal{X}} \left| f(\cdot) = \sum_{i=1}^{\infty} \beta_i k(\cdot, z_i), \beta_i \in \mathbf{R}, z_i \in \mathcal{X}, \|f\| < \infty \right. \right\}. \quad (13)$$

Here, $\|\cdot\|$ is the norm in the RKHS H_k associated with k , i.e. for any $z_i \in \mathcal{X}$, $\beta_i \in \mathbf{R}$ ($i \in \mathbf{N}$),

$$\left\| \sum_{i=1}^{\infty} \beta_i k(\cdot, z_i) \right\|^2 = \sum_{i=1}^{\infty} \sum_{j=1}^{\infty} \beta_i \beta_j k(z_i, z_j). \quad (14)$$

Then any $f \in \mathcal{F}$ minimizing the regularized risk functional

$$c((x_1, y_1, f(x_1)), \dots, (x_m, y_m, f(x_m))) + g(\|f\|) \quad (15)$$

admits a representation of the form

$$f(\cdot) = \sum_{i=1}^m \alpha_i k(\cdot, x_i). \quad (16)$$

Let us give a few remarks before the proof. In its original form, with mean squared loss

$$c((x_1, y_1, f(x_1)), \dots, (x_m, y_m, f(x_m))) = \frac{1}{m} \sum_{i=1}^m (y_i - f(x_i))^2, \quad (17)$$

or hard constraints on the outputs, and $g(\|f\|) = \lambda \|f\|^2$ ($\lambda > 0$), the theorem is due to [15]. Note that in our formulation, hard constraints on the solution are included by the possibility of c taking the value ∞ . A generalization to non-quadratic cost functions was stated by [7], cf. the discussion in [25] (note, however, that [7] did not yet allow for coupling of losses at different points). The present generalization to $g(\|f\|)$ is, to our knowledge, new. For a machine learning point of view on the representer theorem and a variational proof, cf. [12].

The significance of the theorem is that it shows that a whole range of learning algorithms have solutions that can be expressed as expansions in terms of the training examples. Note that monotonicity of g is necessary to ensure that the theorem holds. It does not ensure that the regularized risk functional (15) does not have multiple local minima. For this, we would need to require convexity of g and of the cost function c . If we discarded the strictness of the monotonicity of g , it would no longer follow that each minimizer (there might be multiple ones) of the regularized risk admits an expansion (16); however, it would still follow that there is always another solution minimizing (15) that *does* admit the expansion. In the SV community, (16) is called the *SV expansion*.

Proof. As we have assumed that k maps into \mathbf{R} , we will use (cf. (7))

$$\phi : \mathcal{X} \rightarrow \mathbf{R}^{\mathcal{X}}, \quad x \mapsto k(\cdot, x). \quad (18)$$

Since k is a reproducing kernel, evaluation of the function $\phi(x)$ on the point x' yields

$$(\phi(x))(x') = k(x', x) = \langle \phi(x'), \phi(x) \rangle \quad (19)$$

for all $x, x' \in \mathcal{X}$. Here, $\langle \cdot, \cdot \rangle$ denotes the dot product of H_k .

Given x_1, \dots, x_m , any $f \in \mathcal{F}$ can be decomposed into a part that lives in the span of the $\phi(x_i)$ and a part which is orthogonal to it, i.e.

$$f = \sum_{i=1}^m \alpha_i \phi(x_i) + v \quad (20)$$

for some $\alpha \in \mathbf{R}^m$ and $v \in \mathcal{F}$ satisfying, for all j ,

$$\langle v, \phi(x_j) \rangle = 0. \quad (21)$$

Using the latter and (19), application of f to an arbitrary training point x_j yields

$$f(x_j) = \left\langle \sum_{i=1}^m \alpha_i \phi(x_i) + v, \phi(x_j) \right\rangle = \sum_{i=1}^m \alpha_i \langle \phi(x_i), \phi(x_j) \rangle, \quad (22)$$

independent of v . Consequently, the first term of (15) is independent of v . As for the second term, since v is orthogonal to $\sum_{i=1}^m \alpha_i \phi(x_i)$, and g is strictly monotonic, we get

$$\begin{aligned} g(\|f\|) &= g\left(\left\|\sum_i \alpha_i \phi(x_i) + v\right\|\right) = g\left(\sqrt{\left\|\sum_i \alpha_i \phi(x_i)\right\|^2 + \|v\|^2}\right) \\ &\geq g\left(\left\|\sum_i \alpha_i \phi(x_i)\right\|\right), \end{aligned} \quad (23)$$

with equality occurring if and only if $v = 0$. Setting $v = 0$ thus does not affect the first term of (15), while strictly reducing the second term — hence, any minimizer must have $v = 0$. Consequently, any solution takes the form $f = \sum_i \alpha_i \phi(x_i)$, i.e., using (19),

$$f(\cdot) = \sum_{i=1}^m \alpha_i k(\cdot, x_i). \quad (24)$$

The theorem is proven.

The extension to the case where we also include a parametric part is slightly more technical but straightforward. We state the corresponding result without proof:

Theorem 2 (Semiparametric Representer Theorem). *Suppose that in addition to the assumptions of the previous theorem we are given a set of M real-valued functions $\{\psi_p\}_{p=1}^M$ on \mathcal{X} , with the property that the $m \times M$ matrix $(\psi_p(x_i))_{ip}$ has rank M . Then any $\tilde{f} := f + h$, with $f \in \mathcal{F}$ and $h \in \text{span}\{\psi_p\}$, minimizing the regularized risk*

$$c\left((x_1, y_1, \tilde{f}(x_1)), \dots, (x_m, y_m, \tilde{f}(x_m))\right) + g(\|f\|) \quad (25)$$

admits a representation of the form

$$\tilde{f}(\cdot) = \sum_{i=1}^m \alpha_i k(x_i, \cdot) + \sum_{p=1}^M \beta_p \psi_p(\cdot), \quad (26)$$

with unique coefficients $\beta_p \in \mathbf{R}$ for all $p = 1, \dots, M$.

Remark 1 (Biased regularization). A straightforward extension of the representer theorems can be obtained by including a term $-\langle f_0, f \rangle$ into (15) or (25), respectively, where $f_0 \in H_k$. In this case, if a solution to the minimization problem exists, it admits an expansion which differs from the above ones in that it additionally contains a multiple of f_0 . To see this, decompose the vector v used in the proof of Theorem 1 into a part orthogonal to f_0 and the remainder.

In the case where $g(\|f\|) = \frac{1}{2}\|f\|^2$, this can be seen to correspond to an effective overall regularizer of the form $\frac{1}{2}\|f - f_0\|^2$. Thus, it is no longer the size of $\|f\|$ that is penalized, but the difference to f_0 .

Some explicit applications of Theorems 1 and 2 are given below.

Example 1 (SV regression). For SV regression with the ε -insensitive loss [23] we have

$$c\left((x_i, y_i, f(x_i))_{i=1, \dots, m}\right) = \frac{1}{\lambda} \sum_{i=1}^m \max(0, |f(x_i) - y_i| - \varepsilon) \quad (27)$$

and $g(\|f\|) = \|f\|^2$, where $\lambda > 0$ and $\varepsilon \geq 0$ are fixed parameters which determine the trade-off between regularization and fit to the training set. In addition, a single ($M = 1$) constant function $\psi_1(x) = b$ ($b \in \mathbf{R}$) is used as an offset that is not regularized by the algorithm [25].

In [22], a semiparametric extension was proposed which shows how to deal with the case $M > 1$ algorithmically. Theorem 2 applies in that case, too.

Example 2 (SV classification). Here, the targets satisfy $y_i \in \{\pm 1\}$, and we use

$$c\left((x_i, y_i, f(x_i))_{i=1, \dots, m}\right) = \frac{1}{\lambda} \sum_{i=1}^m \max(0, 1 - y_i f(x_i)), \quad (28)$$

the regularizer $g(\|f\|) = \|f\|^2$, and $\psi_1(x) = b$. For $\lambda \rightarrow 0$, we recover the hard margin SVM, i.e., the minimizer must correctly classify each training point (x_i, y_i) . Note that after training, the actual classifier will be $\text{sgn}(f(\cdot) + b)$.

Example 3 (SVMs minimizing actual risk bounds). The reason why SVM algorithms such as the ones discussed above use the regularizer $g(\|f\|) = \|f\|^2$ are practical ones. It is usually argued that theoretically, we should really be minimizing an upper bound on the expected test error, but in practice, we use a quadratic regularizer, traded off with an empirical risk term via some constant. One can show that in combination with certain loss functions (hard constraints, linear loss, quadratic loss, or suitable combinations thereof), this regularizer

leads to a convex quadratic programming problem [5,23]. In that case, the standard Kuhn-Tucker machinery of optimization theory [19] can be applied to derive a so-called dual optimization problem, which consists of finding the expansion coefficients $\alpha_1, \dots, \alpha_m$ rather than the solution f in the RKHS.

From the point of view of learning theory, on the other hand, more general functions g might be preferable, such as ones that are borrowed from uniform convergence bounds. For instance, we could take inspiration from the basic pattern recognition bounds of [23] and use, for some small $\eta > 0$, the (strictly monotonic) function

$$g(\|f\|) := \sqrt{\frac{R^2\|f\|^2 \left(\log \frac{2m}{R^2\|f\|^2} + 1 \right) - \log(\eta/4)}{m}}. \quad (29)$$

More sophisticated functions based on other bounds on the test error are conceivable, as well as variants for regression estimation (e.g., [3,6]).

Unlike Wahba's original version, the generalized representer theorem can help dealing with these cases. It asserts that the solution still has an expansion in terms of the training examples. It is thus sufficient to minimize the risk bound over expansions of the form (16) (or (26)). Substituting (16) into (15), we get an (m -dimensional) problem in coefficient representation (without having to appeal to methods of optimization theory)

$$\begin{aligned} \min_{\alpha_1, \dots, \alpha_m \in \mathbf{R}} \quad & c \left((x_1, y_1, \sum_{i=1}^m \alpha_i k(x_1, x_i)), \dots, (x_m, y_m, \sum_{i=1}^m \alpha_i k(x_m, x_i)) \right) \\ & + g \left(\sqrt{\sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j k(x_i, x_j)} \right). \end{aligned} \quad (30)$$

If g and c are convex, then so will be the dual, thus we can solve it employing methods of convex optimization (such as interior point approaches often used even for standard SVMs). If the dual is non-convex, we will typically only be able to find local minima.

Independent of the convexity issue, the result lends itself well to gradient-based on-line algorithms for minimizing RKHS-based risk functionals [10,9,17, 11,8,16]: for the computation of gradients, we only need the objective function to be differentiable; convexity is not required. Such algorithms can thus be adapted to deal with more general regularizers.

Example 4 (Bayesian MAP estimates). The well-known correspondence to Bayesian methods is established by identifying (15) with the negative log posterior [14,12]. In this case, $\exp(-c((x_i, y_i, f(x_i))_{i=1, \dots, m}))$ is the likelihood of the data, while $\exp(-g(\|f\|))$ is the prior over the set of functions. The well-known Gaussian process prior (e.g. [24,27]), with covariance function k , is obtained by using $g(\|f\|) = \lambda\|f\|^2$ (here, $\lambda > 0$, and, as above, $\|\cdot\|$ is the norm of the RKHS associated with k). A Laplacian prior would be obtained by using

$g(\|f\|) = \lambda\|f\|$. In all cases, the minimizer of (15) corresponds to a function with maximal a posteriori probability (MAP).

Example 5 (Kernel PCA). PCA in a kernel feature space can be shown to correspond to the case of

$$c((x_i, y_i, f(x_i))_{i=1, \dots, m}) = \begin{cases} 0 & \text{if } \frac{1}{m} \sum_{i=1}^m \left(f(x_i) - \frac{1}{m} \sum_{j=1}^m f(x_j) \right)^2 = 1 \\ \infty & \text{otherwise} \end{cases} \quad (31)$$

with g an arbitrary strictly monotonically increasing function [21]. The constraint ensures that we are only considering linear feature extraction functionals that produce outputs of unit empirical variance. Note that in this case of unsupervised learning, there are no labels y_i to consider.

3 Conclusion

We have shown that for a large class of algorithms minimizing a sum of an empirical risk term and a regularization term in a reproducing kernel Hilbert space, the optimal solutions can be written as kernel expansions in terms of training examples. This has been known for specific algorithms; e.g., for the SV algorithm, where it is a direct consequence of the structure of the optimization problem, but not in more complex cases, such as the direct minimization of some bounds on the test error (cf. Example 3). The representer theorem puts these individual findings into a wider perspective, and we hope that the reader will find the present generalization useful by either gaining some insight, or by taking it as a practical guideline for designing novel kernel algorithms: as long as the objective function can be cast into the form considered in the generalized representer theorem, one can recklessly carry out algorithms in infinite dimensional spaces, since the solution will always live in a specific subspace whose dimensionality equals at most the number of training examples.

Acknowledgements. Thanks to Bob Williamson, Grace Wahba, Jonathan Baxter, Peter Bartlett, and Nello Cristianini for useful comments. This work was supported by the Australian Research Council. AS was supported by DFG grant SM 62/1-1.

References

1. M. A. Aizerman, É. M. Braverman, and L. I. Rozonoér. Theoretical foundations of the potential function method in pattern recognition learning. *Automation and Remote Control*, 25:821–837, 1964.
2. N. Aronszajn. Theory of reproducing kernels. *Transactions of the American Mathematical Society*, 68:337–404, 1950.

3. P. L. Bartlett and J. Shawe-Taylor. Generalization performance of support vector machines and other pattern classifiers. In B. Schölkopf, C. J. C. Burges, and A. J. Smola, editors, *Advances in Kernel Methods — Support Vector Learning*, pages 43–54, Cambridge, MA, 1999. MIT Press.
4. C. Berg, J. P. R. Christensen, and P. Ressel. *Harmonic Analysis on Semigroups*. Springer-Verlag, New York, 1984.
5. B. E. Boser, I. M. Guyon, and V. N. Vapnik. A training algorithm for optimal margin classifiers. In D. Haussler, editor, *Proceedings of the 5th Annual ACM Workshop on Computational Learning Theory*, pages 144–152, Pittsburgh, PA, July 1992. ACM Press.
6. O. Bousquet and A. Elisseeff. Algorithmic stability and generalization performance. In T. K. Leen, T. G. Dietterich, and V. Tresp, editors, *Advances in Neural Information Processing Systems 13*. MIT Press, 2001.
7. D. Cox and F. O’Sullivan. Asymptotic analysis of penalized likelihood and related estimators. *Annals of Statistics*, 18:1676–1695, 1990.
8. L. Csató and M. Opper. Sparse representation for Gaussian process models. In T. K. Leen, T. G. Dietterich, and V. Tresp, editors, *Advances in Neural Information Processing Systems 13*. MIT Press, 2001.
9. Y. Freund and R. E. Schapire. Large margin classification using the perceptron algorithm. In J. Shavlik, editor, *Machine Learning: Proceedings of the Fifteenth International Conference*, San Francisco, CA, 1998. Morgan Kaufmann.
10. T.-T. Frieß, N. Cristianini, and C. Campbell. The kernel adatron algorithm: A fast and simple learning procedure for support vector machines. In J. Shavlik, editor, *15th International Conf. Machine Learning*, pages 188–196. Morgan Kaufmann Publishers, 1998.
11. C. Gentile. Approximate maximal margin classification with respect to an arbitrary norm. Unpublished.
12. F. Girosi, M. Jones, and T. Poggio. Regularization theory and neural networks architectures. *Neural Computation*, 7(2):219–269, 1995.
13. D. Haussler. Convolutional kernels on discrete structures. Technical Report UCSC-CRL-99-10, Computer Science Department, University of California at Santa Cruz, 1999.
14. G. S. Kimeldorf and G. Wahba. A correspondence between Bayesian estimation on stochastic processes and smoothing by splines. *Annals of Mathematical Statistics*, 41:495–502, 1970.
15. G. S. Kimeldorf and G. Wahba. Some results on Tchebycheffian spline functions. *J. Math. Anal. Appl.*, 33:82–95, 1971.
16. J. Kivinen, A. J. Smola, P. Wankadia, and R. C. Williamson. On-line algorithms for kernel methods. in preparation, 2001.
17. A. Kowalczyk. Maximal margin perceptron. In A. J. Smola, P. L. Bartlett, B. Schölkopf, and D. Schuurmans, editors, *Advances in Large Margin Classifiers*, pages 75–113, Cambridge, MA, 2000. MIT Press.
18. H. Lodhi, J. Shawe-Taylor, N. Cristianini, and C. Watkins. Text classification using string kernels. Technical Report 2000-79, NeuroCOLT, 2000. Published in: T. K. Leen, T. G. Dietterich and V. Tresp (eds.), *Advances in Neural Information Processing Systems 13*, MIT Press, 2001.
19. O. L. Mangasarian. *Nonlinear Programming*. McGraw-Hill, New York, NY, 1969.
20. J. Mercer. Functions of positive and negative type and their connection with the theory of integral equations. *Philosophical Transactions of the Royal Society, London*, A 209:415–446, 1909.

21. B. Schölkopf, A. Smola, and K.-R. Müller. Kernel principal component analysis. In B. Schölkopf, C. J. C. Burges, and A. J. Smola, editors, *Advances in Kernel Methods - Support Vector Learning*, pages 327–352. MIT Press, Cambridge, MA, 1999.
22. A. Smola, T. Frieß, and B. Schölkopf. Semiparametric support vector and linear programming machines. In M. S. Kearns, S. A. Solla, and D. A. Cohn, editors, *Advances in Neural Information Processing Systems 11*, pages 585–591, Cambridge, MA, 1999. MIT Press.
23. V. Vapnik. *The Nature of Statistical Learning Theory*. Springer, NY, 1995.
24. G. Wahba. *Spline Models for Observational Data*, volume 59 of *CBMS-NSF Regional Conference Series in Applied Mathematics*. SIAM, Philadelphia, 1990.
25. G. Wahba. Support vector machines, reproducing kernel Hilbert spaces and the randomized GACV. In B. Schölkopf, C. J. C. Burges, and A. J. Smola, editors, *Advances in Kernel Methods — Support Vector Learning*, pages 69–88, Cambridge, MA, 1999. MIT Press.
26. C. Watkins. Dynamic alignment kernels. In A. J. Smola, P. L. Bartlett, B. Schölkopf, and D. Schuurmans, editors, *Advances in Large Margin Classifiers*, pages 39–50, Cambridge, MA, 2000. MIT Press.
27. C. K. I. Williams. Prediction with Gaussian processes: From linear regression to linear prediction and beyond. In M. I. Jordan, editor, *Learning and Inference in Graphical Models*. Kluwer, 1998.